3.1

Def$^n$: a is a subgradient of f at $x$ if
$$f(y) \geq f(x) + a^T(y-x) \quad \forall y$$

* Similarly, since the dual function g is concave, a is a supergradient of g at $x$ if
$$g(y) \leq g(x) + a^T(y-x) \quad \forall y.$$

minimize $f(w)$
s.t $Aw = b$

Then; $g(\lambda) = \inf_w L(w, \lambda) = f(w) + \lambda^T(Aw - b)$

Then, for $\lambda_0 \in \mathbb{R}^n$,
$$g(\lambda_0) = \inf_w f(w) + \lambda_0^T(Aw - b)$$

Let $w_{\lambda_0} = \operatorname{argmin}_w f(w) + \lambda_0^T(Aw - b)$

$\implies g(\lambda_0) = f(w_{\lambda_0}) + \lambda_0^T(Aw_{\lambda_0} - b)$ ——①

Similarly, ~~let~~ for any $\lambda \in \mathbb{R}^n$ let
$$w_\lambda = \operatorname{argmin}_w ~~f(w)~~ f(w) + \lambda^T(Aw - b)$$

$\implies g(\lambda) = f(w_\lambda) + \lambda^T(Aw_\lambda - b)$ ——②

Note that

$$g(\lambda) = \inf_{w} \; f(w) + \lambda^T (Aw - b)$$

$$\leq f(w) + \lambda^T (Aw - b) \quad \text{for all } w \quad \boxed{3}$$

$\therefore$ ③ is true for $w = w_{\lambda_0}$

$$\implies g(\lambda) \leq f(w_{\lambda_0}) + \lambda^T (Aw_{\lambda_0} - b) \quad \text{———} \quad \boxed{4}$$

Then, substituting $f(w_{\lambda_0})$ using ① in ④ gives

$$g(\lambda) \leq g(\lambda_0) - \lambda_0^T (Aw_{\lambda_0} - b) + \lambda^T (Aw_{\lambda_0} - b)$$

$$= g(\lambda_0) + (\lambda - \lambda_0)^T (Aw_{\lambda_0} - b)$$

$$= g(\lambda_0) + (Aw_{\lambda_0} - b)^T (\lambda - \lambda_0)$$

$$\implies \left( Aw_{\lambda_0} - b \right) \text{ is a supergradient of } g$$

at $w_{\lambda_0}$

$$\therefore \; (Aw - b) \in \partial g(\lambda)$$

$$\uparrow \text{ set of supergradients of } g \text{ at } \lambda$$



Group 3: Nice work!

~~Since f is smooth and strongly convex, g is also smooth and strongly convex~~

Since f is $L$-smooth and $\mu$-strongly convex, g is also smooth and strongly con~~vex~~cave with constants

$$L_g = \frac{\lambda_{max}(AA^T)}{\mu} \quad \text{and}$$

$$\mu_g = \frac{\lambda_{min}^+(AA^T)}{L}, \quad \text{respectively}, \quad \text{where}$$

$\lambda_{max}$ is the maximum eigenvalue of $AA^T$ and $\lambda_{min}^+$ is the minimum nonzero eigenvalue of $AA^T$

(Above is a well-known theorem)

* Let $\ell(\lambda) = -g(\lambda)$ (the negative dual fu$^n$)

Then $\nabla \ell(\lambda) = -\nabla g(\lambda)$

since $\ell$ is $L_g$ smooth, we have

$$\ell(\lambda_2) \leq \ell(\lambda_1) + \nabla\ell(\lambda_1)^T(\lambda_2 - \lambda_1) + \frac{L_g}{2}\|\lambda_2 - \lambda_1\|^2 \qquad ①$$

Let $\lambda_2 = \lambda^{k+1}$ & $\lambda_1 = \lambda^k$. Then

$$\lambda_2 - \lambda_1 = \lambda^{k+1} - \lambda^k = -\alpha_k \nabla\ell(\lambda^k),$$

$$(\because \lambda^{k+1} = \lambda^k - \alpha_k\nabla\ell(\lambda^k))$$

∴ ① ⟹

$$\ell(\lambda^{k+1}) \leq \ell(\lambda^k) + \nabla\ell(\lambda^k)^T(-\alpha_k\nabla\ell(\lambda^k))$$
$$+ \frac{L_g}{2}\alpha_k^2\|\nabla\ell(\lambda^k)\|^2$$

$$\Rightarrow \ell(a^{k+1}) \le \ell(a^k) - \alpha_k \|\ell(a^k)\|^2 + \frac{L_g}{2} \alpha_k^2 \|\nabla \ell(a^k)\|^2$$

$$\le \ell(a^k) - \alpha_k \|\ell(a^k)\|^2$$
$$+ \frac{L_g}{2} \alpha_k \times \frac{1}{L_g} \|\nabla \ell(a^k)\|^2$$

for all $\quad \alpha_k \le \frac{1}{L_g}$

$$= \ell(a^k) - \frac{\alpha_k}{2} \|\nabla \ell(a^k)\|^2 \quad \text{———②}$$

Since $\ell(a)$ is strongly convex with constant $m_g$, we have that

$$\|\nabla h(a^k)\|^2 \ge 2 m_g \left( \ell(a^k) - \ell(a^*) \right)$$

Then ② $\Rightarrow$

$$\ell(a^{k+1}) \le \ell(a^k) - \frac{\alpha_k}{2} \times 2 m_g \left( \ell(a^k) - \ell(a^*) \right)$$

$$= \left( 1 - \alpha_k m_g \right) \ell(a^k) + \alpha_k m_g \ell(a^*)$$

$$\Rightarrow \ell(a^{k+1}) - \ell(a^*) \le \left( 1 - \alpha_k m_g \right) \left( \ell(a^k) - \ell(a^*) \right) \quad \text{③}$$

Let $\alpha_k = \alpha$ (constant stem size).

Then by using recursive application of ③ we get

$$\ell(a^k) - \ell(a^*) \le \left( 1 - \alpha m_g \right)^k \left[ \ell(a^0) - \ell(a^*) \right] \quad \text{④}$$

~~Here 1 - αrg~~

Here $0 < 1 - \alpha r_g < 1$, since $0 < \alpha < \frac{1}{L_g}$
and $r_g \leq L_g$.

$\therefore$ ④ $\implies$ $h(a^k)$ converges to the optimal value $h(a^*)$ with a linear rate.

Group 3: Nice!

* The solution $w_{1k+1} = \arg\min_{w}(L w, \lambda_k)$

is primal freasible only when the algorithm converges to it's optimal value. ~~In~~ In the intermediate itterations, the solution $w_{k+1}$ is not ~~prima~~ feasible.

Group 3: Would be great if you can elaborate this a little bit.

3.3

$(P_2):$ minimize $\frac{1}{N} \sum f_i(w_i)$
s.t $w_i = w_j$ $\forall \ j \in N_i$ , $i \in [N]$

Write $P_2$ equivalently as

minimize $\frac{1}{N} \sum_{i \in [N]} f_i(w_i)$

s.t $a_{ij}(w_i - w_j) = 0$ $\forall \ j \in N_i, \ i \in [N]$,

where $A = [a_{ij}]$ is a doubly stochastic matrix.

Write $(P_2)$ equivalently as

$$\text{minimize} \quad \frac{1}{N} \sum_{i \in [N]} f_i(w_i)$$

$$\text{s.t} \quad a_{ij}(w_i - w_j) = 0 \quad \forall \; j \in N_i \quad , \text{where}$$

$A = [a_{ij}]$ is a doubly stochastic matrix compatible with an arbitrary undirected and connected graph.

Then adding all the constraints associated with each $j \in N_i$ gives us

$$\sum_{j \in N_i} a_{ij}(w_i - w_j) = 0$$

$$\implies \sum_{j \in N_i} a_{ij} w_i = \sum_{j \in N_i} a_{ij} w_j$$

$$\implies \left( \sum_{j \in N_i} a_{ij} \right) w_i = \sum_{j \in N_i} a_{ij} w_j$$

$$\underbrace{\qquad}_{=1}$$

$$\implies w_i = \sum_{j \in N_i} a_{ij} w_j \quad \forall \; i \in [N]$$

Then $(P_2)$ can reformulate as

$$\text{minimize}_{w} \quad \frac{1}{N} \sum_{i \in N} f_i(w_i)$$

$$\text{s.t} \quad w_i - \sum_{j \in N_i} a_{ij} w_j = 0 \quad ; \; i \in [N]$$

Then ~~$\partial(\lambda)$~~ the dual function of $(P2)$ is

$$g(\lambda) = \inf_{w} \left[ \sum_{i \in [N]} \frac{f_i(w_i)}{N} + \sum_{i \in [N]} \lambda_i^T \left( w_i - \sum_{j \in N_i} a_{ij} w_j \right) \right]$$

$$= \inf_{w} \left[ \sum_{i \in [N]} \left\{ \frac{f_i(w_i)}{N} + \lambda_i^T \left( w_i - \sum_{j \in N_i} a_{ij} w_j \right) \right\} \right]$$

$$= \sum_{i \in [N]} \inf_{w_i} \left( \frac{f_i(w_i)}{N} + \lambda_i^T \left( w_i - \sum_{j \in N_i} a_{ij} w_j \right) \right)$$

$\underbrace{\phantom{xxxxxxxxxxxx}}$
Subproblems.

$\left( w_i - a_{ii} w_i \right.$
$\left. - \sum_{j \in N_i \setminus \{i\}} a_{ij} w_j \right)$

Then the corresponding ~~dual~~ distributed dual ascent algorithm is

$$w_i^{(k+1)} = \operatorname{argmin}_{w_i} \frac{f_i(w_i)}{N} + \left( \lambda_i^T \right)^{(k)} w_i - \sum a_{ij} w$$

primal
variable $\Rightarrow$ $w_i^{(k+1)} = \operatorname{argmin}_{w_i} \frac{f_i(w_i)}{N} + \left( \lambda_i^T \right)^{(k)} \left[ (1 - a_{ii}) w_i - \sum_{j \in N_i - \{i\}} a_{ij} w_j^{(k)} \right]$
update

dual variable
component $\Rightarrow$ $\lambda_i^{(k+1)} = \lambda_i^{(k)} + \alpha_x \left[ (1 - a_{ii}) w_i^{(k+1)} - \sum_{j \in N_i - \{i\}} a_{ij} w_j \right]$
update

Comparison between the dual method and the primal method (numerically) ~~for a~~ using a particular ~~a~~ connected and undirected graph.

We consider the following communication graph with 5 users.



We used:

* $f_i(W_i) = W_i^T B_i W_i + q_i^T W_i$, ~~where~~ $i = 1, 2, 3, 4, 5$, where $B_i$'s are positive definite matrices.
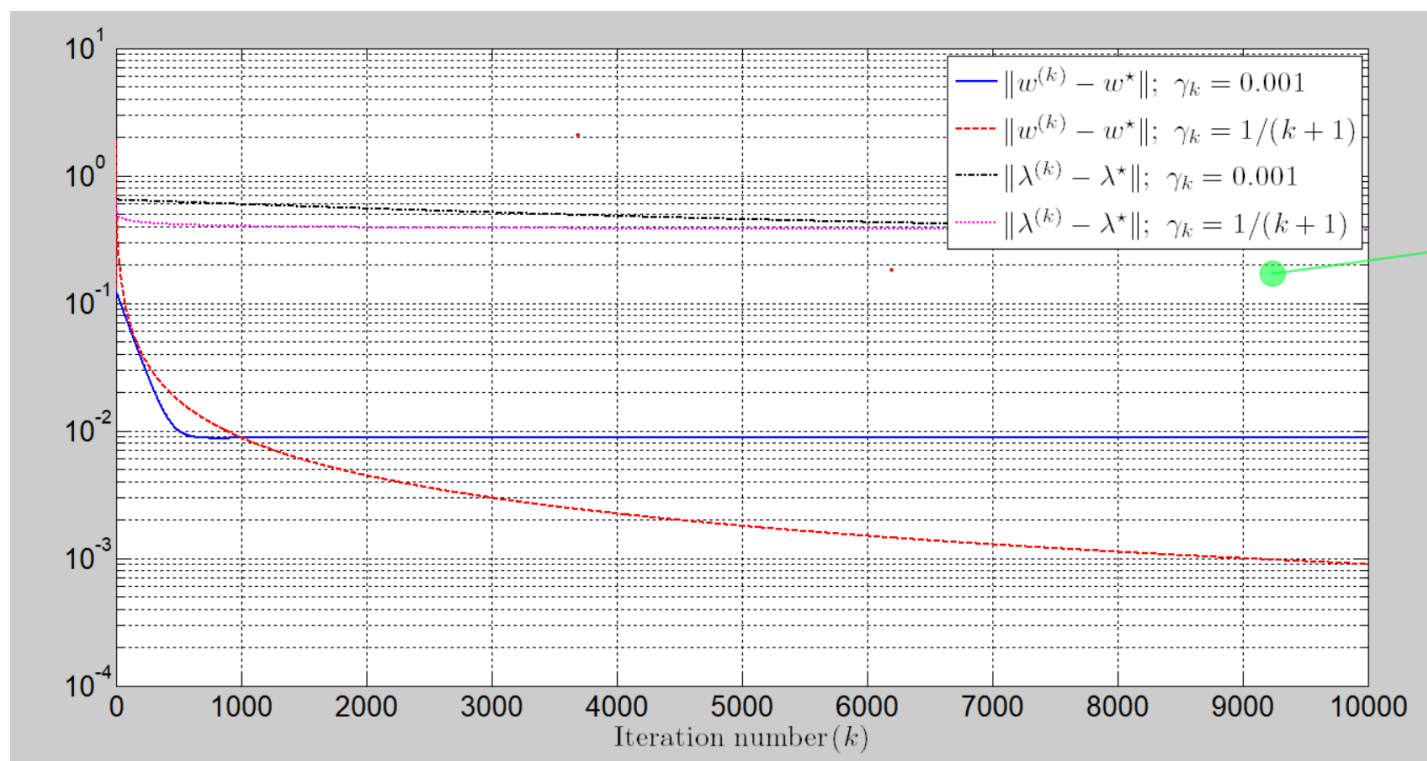
* $W_i \in \mathbb{R}$

* Doubly stochastic matrix $A$ is taken as

$$
A = \begin{bmatrix}
0.75 & 0.25 & 0 & 0 & 0 \\
0.25 & 0.25 & 0.25 & 0.25 & 0 \\
0 & 0.25 & 0.75 & 0 & 0 \\
0 & 0.25 & 0 & \frac{5}{12} & \frac{1}{3} \\
0 & 0 & 0 & \frac{1}{3} & \frac{2}{3}
\end{bmatrix}
$$

Following figure shows the convergences of $\|w^{(t)} - w^*\|$ and $\|\lambda^{(t)} - \lambda^*\|$ using the primal method and the dual method, respectively.

Group 3: Would be great if you could add the code here.

Legend:
- $\|w^{(k)} - w^\star\|$; $\gamma_k = 0.001$
- $\|w^{(k)} - w^\star\|$; $\gamma_k = 1/(k+1)$
- $\|\lambda^{(k)} - \lambda^\star\|$; $\gamma_k = 0.001$
- $\|\lambda^{(k)} - \lambda^\star\|$; $\gamma_k = 1/(k+1)$

Group 3: Would be great if you can add the name of each method in the legend!

According to the graph, for this particular case, the primal method shows better convergences than the dual method. For constant step size, the primal method ~~shows~~ clearly shows a linear rate of convergence. However, with constant stepsize, $w^{(k)}$ seems to converge into a neighbourhood of $w^*$, while with nonsummable and square summable stepsize ~~$(1/k+1)$~~ ($\alpha_k = 1/k+1$) ~~#~~ $w^{(k)}$ converges more towards the optimal solution $w^*$.

In both methods, at every iteration, each ~~the~~ user ~~communication with~~ communicate only with ~~#~~ ~~#~~ it's neighbours. ∴ the communication cost in both methods ~~per iteration~~ are same, per iteration.

Group 3: Good work in overall! :)